

This Page Is Inserted by IFW Operations
and is not a part of the Official Record

BEST AVAILABLE IMAGES

Defective images within this document are accurate representations of the original documents submitted by the applicant.

Defects in the images may include (but are not limited to):

- BLACK BORDERS
- TEXT CUT OFF AT TOP, BOTTOM OR SIDES
- FADED TEXT
- ILLEGIBLE TEXT
- SKEWED/SLANTED IMAGES
- COLORED PHOTOS
- BLACK OR VERY BLACK AND WHITE DARK PHOTOS
- GRAY SCALE DOCUMENTS

IMAGES ARE BEST AVAILABLE COPY.

**As rescanning documents *will not* correct images,
please do not report the images to the
Image Problems Mailbox.**



PATENT ABSTRACTS OF JAPAN

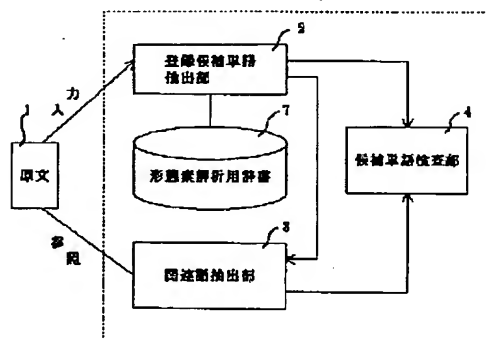
(11) Publication number: **11134334 A**(43) Date of publication of application: **21.05.99**

(51) Int. Cl. **G06F 17/27**
G06F 17/22
G06F 17/28
G06F 17/30

(21) Application number: **09296768**(71) Applicant: **FUJITSU LTD**(22) Date of filing: **29.10.97**(72) Inventor: **SATSUSANO YUKARI****(54) WORD REGISTERING DEVICE AND RECORDING MEDIUM****(57) Abstract:**

PROBLEM TO BE SOLVED: To reduce a labor at the time of a registering work by extracting a word to be registered including an unregistered composite word.

SOLUTION: This device is provided with a morpheme analysis dictionary 7 for processing natural language, and an registration candidate word extracting part 2 for operating the morpheme analysis of a natural language sentence, extracts a composite word in which words which are not registered in the morpheme analysis dictionary 7 and nouns which are not registered in the dictionary 7 are continued, and judges the composite word whose frequency is high as a registration candidate word to be registered.



COPYRIGHT: (C)1999,JPO

(19)日本国特許庁 (J P)

(12) 公 開 特 許 公 報 (A)

(11)特許出願公開番号

特開平11-134334

(43)公開日 平成11年(1999) 5月21日

(51)Int.Cl.⁶

識別記号

F I

G 0 6 F 17/27
17/22
17/28
17/30

G 0 6 F 15/38
15/20
15/38
15/403

E
5 2 2 L
U
3 3 0 C

審査請求 未請求 請求項の数4 O L (全 19 頁)

(21)出願番号

特願平9-296768

(22)出願日

平成9年(1997)10月29日

(71)出願人 000005223

富士通株式会社

神奈川県川崎市中原区上小田中4丁目1番
1号

(72)発明者 堀々野 由香梨

神奈川県川崎市中原区上小田中4丁目1番
1号 富士通株式会社内

(74)代理人 弁理士 山谷 昭榮 (外2名)

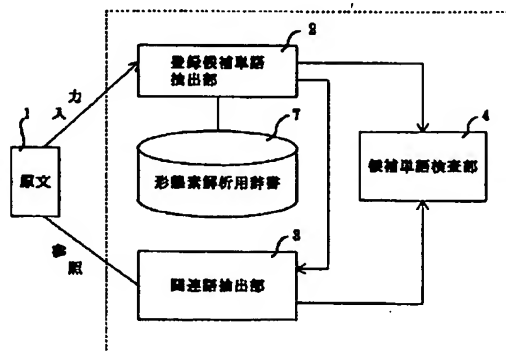
(54)【発明の名称】 単語登録装置及び記録媒体

(57)【要約】

【課題】未登録複合語を含めた登録すべき単語を抽出し、登録作業時の労力を軽減すること。

【解決手段】自然言語を処理するための形態素解析用辞書7と、自然言語文を形態素解析し、前記形態素解析用辞書7に登録されていない単語及び該辞書7に登録されていない名詞類の連続した複合語を抽出して、頻度の高いものを登録すべき登録候補単語と判定する登録候補単語抽出部2とを備える。

本発明の原理説明図



【特許請求の範囲】

【請求項1】自然言語を処理するための形態素解析用辞書と、

自然言語文を形態素解析し、前記形態素解析用辞書に登録されていない単語及び該辞書に登録されていない名詞類の連続した複合語を抽出して、頻度の高いものを登録すべき登録候補単語と判定する登録候補単語抽出部とを備えることを特徴とした単語登録装置。

【請求項2】前記判定した登録候補単語を含む原文を検索し、前記形態素解析用辞書に登録されていない単語及び該辞書に登録されていない名詞類の連続した複合語を抽出する関連語抽出部を備えることを特徴とした請求項1記載の単語登録装置。

【請求項3】前記判定した登録候補単語を含む原文に対して、前記登録候補単語を取り入れる前の形態素解析結果と前記登録候補単語を取り入れた場合の形態素解析結果を比較して、解析誤りが起こっているかどうかを判定する候補単語検査部を備えることを特徴とした請求項1記載の単語登録装置。

【請求項4】コンピュータに、
自然言語文を形態素解析する解析手順と、
前記形態素解析結果から形態素解析用辞書に登録されていない単語を抽出する抽出手順と、
前記形態素解析結果から形態素解析用辞書に登録されていない名詞類の連続した複合語を抽出する抽出手順と、
前記抽出手順で抽出した単語及び複合語より頻度の高い語を登録候補単語と判定する判定手順と、
を実行するためのプログラムを格納したコンピュータ読取可能な記録媒体。

【発明の詳細な説明】

【0001】

【発明の属する技術分野】本発明は、自然言語を処理するための単語辞書に単語を追加登録する単語登録装置及び記録媒体に関する。

【0002】

【従来の技術】日本語を形態素に分割する形態素解析は、自然言語処理の最も基本となる処理である。従来、形態素解析は、自然言語処理の様々なアプリケーションに用いられており、例えば、情報検索や文書中の誤りを発見する文書校正支援に用いられている。

【0003】形態素解析にあたって、それに用いられる形態素解析用辞書は、形態素解析の性能を左右する重要な基本データである。この辞書中に単語が登録されていないと、解析が失敗したり、他の語として誤って解析されてしまう。例えば、事故や事件が起こった場合、関連記事を検索するための新しい単語を入力して検索するというニーズが増大しているが、関連の単語が辞書に入っていない場合、目的とする記事が検索できないという事態が生じる可能性がある。そのため、日々増加している新しい事象を表す単語を収集して、形態素解析用辞書に

追加することが重要である。しかし、新しい事象を表す単語は日々増加しているため、登録すべき単語を収集したり、テストする作業には多くの労力がかかっていた。

【0004】従来、形態素解析用辞書に未登録語を登録する場合、形態素解析手段により入力文の解析を行い、その情報を基に入力文中の未登録語を知らせてユーザに登録を促すことが、特開平3-246673号公報に記載されていた。また、未登録語の出現回数を計算して、使用頻度の多いものから優先的に登録することが、特開昭63-208167号公報に記載されていた。また、既知語の意味カテゴリを用いて未知語の意味カテゴリを推定して登録することが、特開平8-16597号公報に記載されていた。また、関連情報辞書登録手段により、格の違いによる二重登録を排除し辞書量を少なくすることが、特開平6-119374号公報に記載されていた。

【0005】

【発明が解決しようとする課題】前記のような従来のものは、次のような課題があった。

①：二つ以上の名詞類が連続している未登録複合語を抽出できるものではなかった。

②：登録候補単語の関連である入力文中に含まれる頻度の低い未登録語を登録できるものではなかった。

③：登録すべき単語のテストを事前に行えるものではなかった。

【0007】本発明は、このような従来の課題を解決し、未登録複合語を含めた登録すべき単語の抽出をし、登録作業時の労力を軽減し、更に登録すべき単語のテストを事前に行い、質のよい単語を半自動的に収集すること、また、登録すべき単語候補として選ばれた単語と関連のある語も同時に収集できるようにすることを目的とする。

【0008】

【課題を解決するための手段】図1は本発明の原理説明図である。図1中、1は原文、2は登録候補単語抽出部、3は関連語抽出部、4は候補単語検査部、7は形態素解析用辞書である。

【0009】本発明は前記従来の課題を解決するため次のように構成した。

(1)：自然言語を処理するための形態素解析用辞書7と、自然言語文を形態素解析し、前記形態素解析用辞書7に登録されていない単語及び該辞書7に登録されていない名詞類の連続した複合語を抽出して、頻度の高いものを登録すべき登録候補単語と判定する登録候補単語抽出部2とを備える。

【0010】(2)：前記(1)の単語登録装置において、前記判定した登録候補単語を含む原文1を検索し、前記形態素解析用辞書7に登録されていない単語及び該辞書7に登録されていない名詞類の連続した複合語を抽

出する関連語抽出部3を備える。

【0011】(3):前記(1)の単語登録装置において、前記判定した登録候補単語を含む原文1に対して、前記登録候補単語を取り入れる前の形態素解析結果と前記登録候補単語を取り入れた場合の形態素解析結果を比較して、解析誤りが起こっているかどうかを判定する候補単語検査部4を備える。

【0012】(4):コンピュータに、自然言語文を形態素解析する解析手順と、前記形態素解析結果から形態素解析用辞書7に登録されていない単語を抽出する抽出手順と、前記形態素解析結果から形態素解析用辞書7に登録されていない名詞類の連続した複合語を抽出する抽出手順と、前記抽出手順で抽出した単語及び複合語より頻度の高い語を登録候補単語と判定する判定手順と、を実行するためのプログラムを格納したコンピュータ読取可能な記録媒体とする。

【0013】(作用)前記構成に基づく作用を説明する。登録候補単語抽出部2で、自然言語文を形態素解析し、形態素解析用辞書7に登録されていない単語及び該辞書7に登録されていない名詞類の連続した複合語を抽出して、頻度の高いものを登録すべき登録候補単語と判定する。このため、頻度の高い未登録語だけでなく頻度の高い未登録複合語も登録候補単語として判定することができ、登録すべき語の抽出及び選択作業を軽減することができる。

【0014】また、関連語抽出部3で、前記判定した登録候補単語を含む原文1を検索し、形態素解析用辞書7に登録されていない単語及び該辞書7に登録されていない名詞類の連続した複合語を抽出する。このため、頻度が低い語も関連語として原文から抽出し、その語も登録候補単語として取り入れることができる。

【0015】さらに、候補単語検査部4で、前記判定した登録候補単語を含む原文1に対して、前記登録候補単語を取り入れる前の形態素解析結果と前記登録候補単語を取り入れた場合の形態素解析結果を比較して、解析誤りが起こっているかどうかを判定する。このため、登録する前にテストが行え、質のよい単語を収集することができる。

【0016】また、自然言語文を形態素解析する解析手順と、前記形態素解析結果から形態素解析用辞書7に登録されていない単語を抽出する抽出手順と、前記形態素解析結果から形態素解析用辞書7に登録されていない名詞類の連続した複合語を抽出する抽出手順と、前記抽出手順で抽出した単語及び複合語より頻度の高い語を登録候補単語と判定する判定手順と、を実行するためのプログラムを格納したコンピュータ読取可能な記録媒体とする。このため、この記録媒体のプログラムをコンピュータにインストールすることで、頻度の高い未登録語だけでなく頻度の高い未登録複合語も登録候補単語として判定することができる単語登録装置を容易に提供すること

ができる。

【0017】

【発明の実施の形態】本発明の単語登録装置では、日々更新されるニュース記事やWebページ(インターネットのホームページ)等の記事を形態素解析し、登録すべき単語候補を抽出し、その語が登録した場合の解析のテストを行う機構を設けることで、登録すべき単語の抽出や登録作業時の労力を軽減するものである。また、登録すべき単語候補として選ばれた単語と同時に登録すべき関連語も原文から抽出し、その語も登録単語候補として取り入れる機能を備えるものである。

【0018】図2～図16は本発明の実施の形態を示した図である。以下、図2～図16に基づいて本発明の実施の形態を説明する。

(1):装置構成の説明

図2は装置構成図である。図2において、原文データ1が入力される単語登録装置には、登録候補単語抽出部2、関連語抽出部3、候補単語検査部4、単語登録部5、形態素解析エンジン6、形態素解析用辞書7が設けられている。

【0019】原文データ1は、入力手段(図示せず)により入力される日々更新されるニュース記事やWebページ等の記事である。登録候補単語抽出部2は、形態素解析結果から登録候補単語を抽出するものである。関連語抽出部3は、登録候補単語を元に関連語を抽出するものである。候補単語検査部4は、元の解析結果と登録候補単語を取り入れた場合の解析結果を比較して、解析誤りが起こっているかどうかを判定するものである。単語登録部5は、ユーザに登録候補単語や関連語の検査結果を表示し、形態素解析用辞書7に格納するものである。形態素解析エンジン6は、形態素解析を行う処理部である。形態素解析用辞書7は、形態素解析に使用するための単語を登録しておくものである。

【0020】(2):全体の処理手順の説明

図3は全体の処理手順の説明図である。以下、図3の処理S1～処理S4に従って説明する。

【0021】S1:決められた時間にダウンロード等で自動で入力された新聞記事等の原文データ1を登録候補単語抽出部2で、形態素解析し、その結果から登録候補単語を抽出し、処理S2に移る。

【0022】S2:関連語抽出部3で、登録候補単語として選ばれた単語を含む元記事中に含む単語(関連語)を登録候補単語として選択し、処理S3に移る。

S3:候補単語検査部4で、登録候補単語及び関連語を登録した場合の形態素解析結果をテストし、その結果をユーザに提示し、処理S4に移る。

【0023】S4:ユーザが登録すべき単語として指示した場合、単語登録部5で形態素解析用辞書7に登録して、この処理を終了する。

(3):登録候補単語抽出部の処理の説明

図4は登録候補単語抽出部の処理の説明図である。以下、図4の処理S11～処理S14に従って説明する。

【0024】S11：登録候補単語抽出部2は、原文データ1に対して、形態素解析エンジン6と形態素解析用辞書7を用いて形態素解析を行い、処理S12に移る。

S12：登録候補単語抽出部2は、形態素解析結果から未登録語を抽出して、未登録語頻度表を作成し、処理S13に移る。

【0025】S13：登録候補単語抽出部2は、形態素解析結果から名詞類の連続を抽出して、未登録複合語頻度表を作成し、処理S14に移る。

S14：登録候補単語抽出部2は、それぞれ作成した頻度表の頻度の上位のものを登録候補単語リストに登録して、この処理を終了する。

【0026】(4)：関連語抽出部の処理の説明

図5は関連語抽出部の処理の説明図である。以下、図5の処理S21～処理S23に従って説明する。

【0027】S21：関連語抽出部3は、登録候補単語を含む元の文の記事を検索し、処理S22に移る。

S22：関連語抽出部3は、その記事中に未登録語頻度表、未登録複合語頻度表に含まれる語が存在するかを判定し、処理S23に移る。

【0028】S23：関連語抽出部3は、各頻度表に含まれる語があれば、それを関連語として抽出し、登録候補単語リストに追加して、この処理を終了する。

(5)：候補単語検査部の処理の説明

図6は候補単語検査部の処理の説明図である。以下、図6の処理S31～処理S34に従って説明する。

【0029】S31：候補単語検査部4は、登録候補単語リストから候補単語辞書を作成し、処理S32に移る。

S32：候補単語検査部4は、登録候補単語を含む原文に対して、元の形態素解析用辞書と候補単語辞書を用いて、形態素解析をし、処理S33に移る。

【0030】S33：候補単語検査部4は、元の形態素解析結果と登録候補単語を取り入れた場合の形態素解析結果を比較して、解析誤りが起こっているかどうかを判定し、処理S34に移る。

【0031】S34：候補単語検査部4は、解析誤りが起こっている単語を登録候補単語リストから除外し、この処理を終了する。なお、解析誤りの例として、登録候補単語を取り入れた場合に他の部分（特に取り入れた登録候補単語の前後部分）が未登録語となる場合や逆に未登録語が増加する場合がある。

【0032】(6)：単語登録部の処理の説明

図7は単語登録部の処理の説明図である。以下、図7の処理S41～処理S44に従って説明する。

【0033】S41：単語登録部5は、登録候補単語リストと元の形態素解析結果とそれに新たに登録した場合の形態素解析結果をユーザに提示し、処理S42に移

る。

S42：ユーザが登録候補単語から登録すべき単語を選択し、処理S43に移る。

【0034】S43：単語登録部5は、ユーザに単語の辞書上の登録情報を候補単語辞書から提示し、処理S44に移る。

S44：ユーザが候補単語辞書の内容をそのまま、あるいは修正して、単語登録部5で形態素解析用辞書7に登録し、この処理を終了する。

【0035】(7)：具体例による説明

a：登録候補単語を登録する場合の説明

図8は登録候補単語を登録する場合の説明図(1)であり、図8(a)は一文の形態素解析例の説明、図8

(b)は未登録単語頻度表の説明である。図9は登録候補単語を登録する場合の説明図(2)であり、図9

(a)は候補単語辞書の説明、図9(b)は登録前の形態素解析結果の説明である。図10は登録候補単語を登録する場合の説明図(3)であり、図10(a)は「ヤンゴン」を登録した場合の形態素解析結果の説明、図10(b)はユーザが修正した候補単語辞書の説明である。

【0036】以下は、いくつかの内容を含む新聞記事から登録単語を抽出する例を図8～図10により説明する。まず、登録候補単語抽出部2において、原文を形態素解析する。形態素解析の結果は、例えば、図8(a)のように、文が形態素単位に分割され、それぞれの品詞、詳細品詞、表記が出力される。

【0037】登録候補単語抽出部2では、形態素解析の解析結果から、詳細品詞が「未登録語」となっている単語を収集し、図8(b)のように頻度が記入された未登録単語頻度表を作成する。

【0038】登録候補単語抽出部2は、原文の数に応じて頻度が上位であるものを登録すべき候補の単語として抽出する。例えば、ここで頻度が「10」で頻度の高い「ヤンゴン」を登録候補単語として抽出する。候補単語検査部4では、登録候補単語である「ヤンゴン」に仮の品詞として、普通名詞を付与し、候補単語辞書を作成する。この候補単語辞書は、図9(a)のように表記、品詞、詳細品詞が設けられている候補単語検査部4では、登録候補単語が出現している文を元の形態素解析用辞書7と登録候補単語を取り入れた辞書を使って解析し直して、その結果を出力する。例えば、登録候補単語「ヤンゴン」を含む文が次のものであったとする。

【0039】「ミャンマーの首都ヤンゴンで学生のデモが始まった。」この文に対して、「ヤンゴン」を登録する前の形態素解析結果は、図9(b)であり、「ヤンゴン」を登録した場合の形態素解析結果は、図10(a)である。図9(b)において、未登録語であった「ヤンゴン」は、図10(a)においては普通名詞となり他の単語にも未登録語が含まれていない。このため「ヤンゴ

ン」を登録した場合の結果に解析誤りは含まれていない。

【0040】候補単語検査部4は、この結果を単語登録部5に渡し、ユーザに提示する。ユーザは、この結果を確認し、「ヤンゴン」を辞書に登録することを指示する。ここで「ヤンゴン」は、地名であるので、ユーザは、詳細品詞を「地名」に修正する。即ち、図10

(b)のように候補単語辞書の情報を修正して形態素解析用辞書7に登録する。

【0041】b：未登録複合語頻度表を作成する場合の説明

図11は未登録複合語頻度表を作成する場合の説明図

(1)であり、図11(a)は未登録複合語頻度表の説明、図11(b)は候補単語辞書の説明、図11(c)は登録前の形態素解析結果の説明である。図12は未登録複合語頻度表を作成する場合の説明図(2)であり、図12(a)は登録した後の形態素解析結果の説明、図12(b)はユーザが修正した候補単語辞書の説明である。

【0042】登録候補単語抽出部2で、形態素解析結果から未登録単語頻度表以外に、未登録複合語頻度表を作成するものである。これは、二つ以上の名詞類(名詞、接頭語、接尾語、「・」、「/」、「=」、動詞の連用形等)が連続しているものを取り出し、その頻度を調査したものである。

【0043】ここで、未登録複合語頻度表が、図11(a)のように得られたとする。なお、図11(a)において、形態素の区切りは「/」で表している。ここでは、頻度が「12」と高い「オーム/真理/教」を登録候補単語として抽出したとする。候補単語検査部4では、図11(b)のように「オーム真理教」に仮の品詞として、普通名詞を付与し、この「オーム真理教」が出現した文において形態素解析のテストを行う。

【0044】候補単語検査部4では、登録候補単語が出現している文を元の形態素解析用辞書7と登録候補単語を取り入れた辞書を使って解析し直して、その結果を出力する。ここで、「オーム真理教」を含む原文が次のものであったとする。

【0045】「オーム真理教の信者の林春男容疑者がきょう逮捕されました。」これを「オーム真理教」を一語として登録する前の形態素解析結果は、図11(c)に示してあり、登録した後の形態素解析結果は、図12

(a)に示してある。図11(c)と図12(a)のように、「オーム真理教」を登録した場合の結果に解析誤りは含まれていないので、候補単語検査部4は、この結果を単語登録部5に渡し、ユーザに提示する。

【0046】ユーザは、この結果を確認し、「オーム真理教」を辞書に登録することを指示する。ここで「オーム真理教」は、固有名詞であるので、ユーザは、詳細品詞を「固有名詞」に修正する。即ち、図12(b)のよ

うに候補単語辞書の情報を修正して形態素解析用辞書7に登録する。

【0047】c：関連語を登録する場合の説明

図13は関連語を登録する場合の説明図(1)であり、図13(a)は候補単語辞書(関連語)の説明、図13(b)は登録前の形態素解析結果の説明である。図14は関連語を登録する場合の説明図(2)であり、図14(a)は「國林長」を登録した場合の形態素解析結果の説明、図14(b)は「國林長官狙撃事件」を登録した場合の形態素解析結果の説明である。図15は関連語を登録する場合の説明図(3)であり、図15(a)はユーザが修正した候補単語辞書の説明、図15(b)は登録前の形態素解析結果の説明である。図16は関連語を登録場合の説明図(4)であり、図16(a)は「アウン・タン・スー・チー」を登録した場合の形態素解析結果の説明、図16(b)はユーザが修正した候補単語辞書の説明である。

【0048】前記具体例a、bのように「ヤンゴン」と「オーム真理教」を登録候補単語として抽出した場合、関連語抽出部3では、以下のように処理を行う。関連語抽出部3では、登録候補単語を含む記事中に含まれる頻度の低い未登録語や未登録複合語を選択する。これにより、以下の選択結果が得られたとする。

【0049】「國林長」

「國林長/官/狙撃/事件」

「アウン/・/タン/・/スー/・/チー」

以上の関連語を登録候補単語リストに追加し、候補単語検査部4でテストを行う。候補単語検査部4では、以上の登録候補単語と関連語に仮の品詞として、普通名詞を付与し、それぞれの語が出現した文において形態素解析のテストを行う。例えば、関連語から図13(a)のような候補単語辞書(関連語)を作る。

【0050】候補単語検査部4では、登録候補単語が出現している文を元の形態素解析用辞書7と登録候補単語を取り入れた辞書を使って解析し、その結果を出力する。これは例えば、関連語を含む文が次のようであったとする。

【0051】「警察庁の國林長官狙撃事件の捜査をめぐる対応が適切でない。」

「アウン・タン・スー・チーさんの勢力とは一線を画している。」

・「國林長官狙撃事件」を含む文の形態素解析結果は、登録前は図13(b)となり、「國林長」を登録した場合は図14(a)となり、「國林長官狙撃事件」を登録した場合は図14(b)となる。

【0052】ここで、「國林長」と「國林長官狙撃事件」を登録した場合は、いずれも解析誤りが起こっていないので、候補単語検査部4は、その結果を単語登録部5に渡し、ユーザに提示する。ユーザは、図14(a)と図14(b)の形態素解析結果から、「國林長官狙撃

事件」を登録する方が正しいと判断し、「園林長官狙撃事件」を登録するとユーザが指示する。

【0053】この場合、品詞は固有名詞なので、ユーザは、図15(a)のように候補単語辞書の詳細品詞を「普通名詞」から「固有名詞」に修正し、単語登録部5で形態素解析用辞書7に取り込むようにする。

【0054】次に「アウン・タン・スー・チー」を登録する前と登録した後の形態素解析結果は、図15

(b)と図16(a)のようになる。ここで、「アウン・タン・スー・チー」を登録した場合は、解析誤りが起こっていないので、その結果を単語登録部5に渡し、ユーザに提示する。ユーザは、この結果を確認し、「アウン・タン・スー・チー」を形態素解析用辞書7に登録することを指示する。ここで、「アウン・タン・スー・チー」は人名であるので、ユーザは、候補単語辞書の詳細品詞を「普通名詞」から「人名」に修正し、単語登録部5で形態素解析用辞書7に取り込むようにする。

【0055】以上実施の形態で説明したように、登録すべき単語の抽出および選択が軽減され、更に登録すべき単語のテストを事前に行えるので、質の良い単語を半自動的に収集できる。また、関連のある語も同時に収集することが可能となる。

【0056】(8):プログラムのインストールの説明登録候補単語抽出部2、関連語抽出部3、候補単語検査部4、単語登録部5、形態素解析エンジン6は実際にはプログラムで構成でき、主制御部(CPU)が実行するものであり、主記憶に格納されているものである。これらのプログラムは、一般的な、パーソナルコンピュータ、ワークステーション等のデータ処理装置(コンピュータ)で処理されるものである。これらのコンピュータは、主制御部、主記憶、ハードディスク等のファイル装置、表示装置、キーボード等の入力手段である入力装置などのハードウェアで構成されている。

【0057】このコンピュータに、本発明のプログラムをインストールする。このインストールは、フロッピー、光磁気ディスク等の可搬型の記録媒体に、これらのプログラムを記憶させておき、コンピュータが備えている記憶媒体に対して、アクセスするためのドライブ装置を介して、或いは、LAN等のネットワークを介して、コンピュータに設けられたファイル装置にインストールされる。そして、このファイル装置から処理に必要なプログラムステップを主記憶に読み出し、主制御部が実行するものである。

【0058】

【発明の効果】以上説明したように、本発明によれば次のような効果がある。

(1)登録候補単語抽出部で、自然言語文を形態素解析し、形態素解析用辞書に登録されていない単語及び該辞書に登録されていない名詞類の連続した複合語を抽出して、頻度の高いものを登録すべき登録候補単語と判定す

るため、頻度の高い未登録語だけでなく頻度の高い未登録複合語も登録候補単語として判定することができ、登録すべき語の抽出及び選択作業を軽減することができる。

【0059】(2):関連語抽出部で、登録候補単語を含む原文を検索し、形態素解析用辞書に登録されていない単語及び該辞書に登録されていない名詞類の連続した複合語を抽出するため、頻度が低い単語及び複合語も関連語として原文から抽出し、その語も登録候補単語として取り入れることができる。

【0060】(3):候補単語検査部で、登録候補単語を含む原文に対して、前記登録候補単語を取り入れる前の形態素解析結果と前記登録候補単語を取り入れた場合の形態素解析結果を比較して、解析誤りが起こっているかどうかを判定するため、登録する前にテストが行え、質のよい単語を収集することができる。

【0061】(4):自然言語文を形態素解析する解析手順と、前記形態素解析結果から形態素解析用辞書に登録されていない単語を抽出する抽出手順と、前記形態素解析結果から形態素解析用辞書に登録されていない名詞類の連続した複合語を抽出する抽出手順と、前記抽出手順で抽出した単語及び複合語より頻度の高い語を登録候補単語と判定する判定手順と、を実行するためのプログラムを格納したコンピュータ読取可能な記録媒体とするため、この記録媒体のプログラムをコンピュータにインストールすることで、頻度の高い未登録語だけでなく頻度の高い未登録複合語も登録候補単語として判定することができる単語登録装置を容易に提供することができる。

【図面の簡単な説明】

【図1】本発明の原理説明図である。

【図2】実施の形態における装置構成図である。

【図3】実施の形態における全体の処理手順の説明図である。

【図4】実施の形態における登録候補単語抽出部の処理の説明図である。

【図5】実施の形態における関連語抽出部の処理の説明図である。

【図6】実施の形態における候補単語検査部の処理の説明図である。

【図7】実施の形態における単語登録部の処理の説明図である。

【図8】実施の形態における登録候補単語を登録する場合の説明図(1)である。

【図9】実施の形態における登録候補単語を登録する場合の説明図(2)である。

【図10】実施の形態における登録候補単語を登録する場合の説明図(3)である。

【図11】実施の形態における未登録複合語頻度表を作成する場合の説明図(1)である。

【図 1 2】実施の形態における未登録複合語頻度表を作成する場合の説明図（2）である。

【図 1 3】実施の形態における関連語を登録する場合の説明図（1）である。

【図 1 4】実施の形態における関連語を登録する場合の説明図（2）である。

【図 1 5】実施の形態における関連語を登録する場合の説明図（3）である。

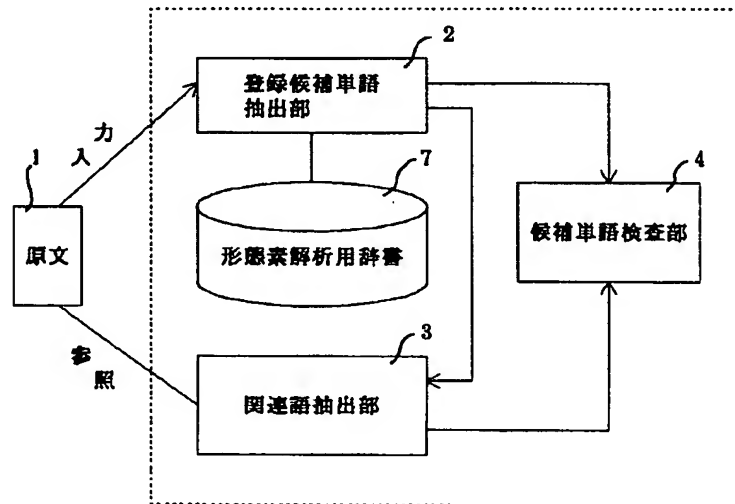
【図 1 6】実施の形態における関連語を登録する場合の説明図（4）である。

【符号の説明】

- 1 原文
- 2 登録候補単語抽出部
- 3 関連語抽出部
- 4 候補単語検査部
- 7 形態素解析用辞書

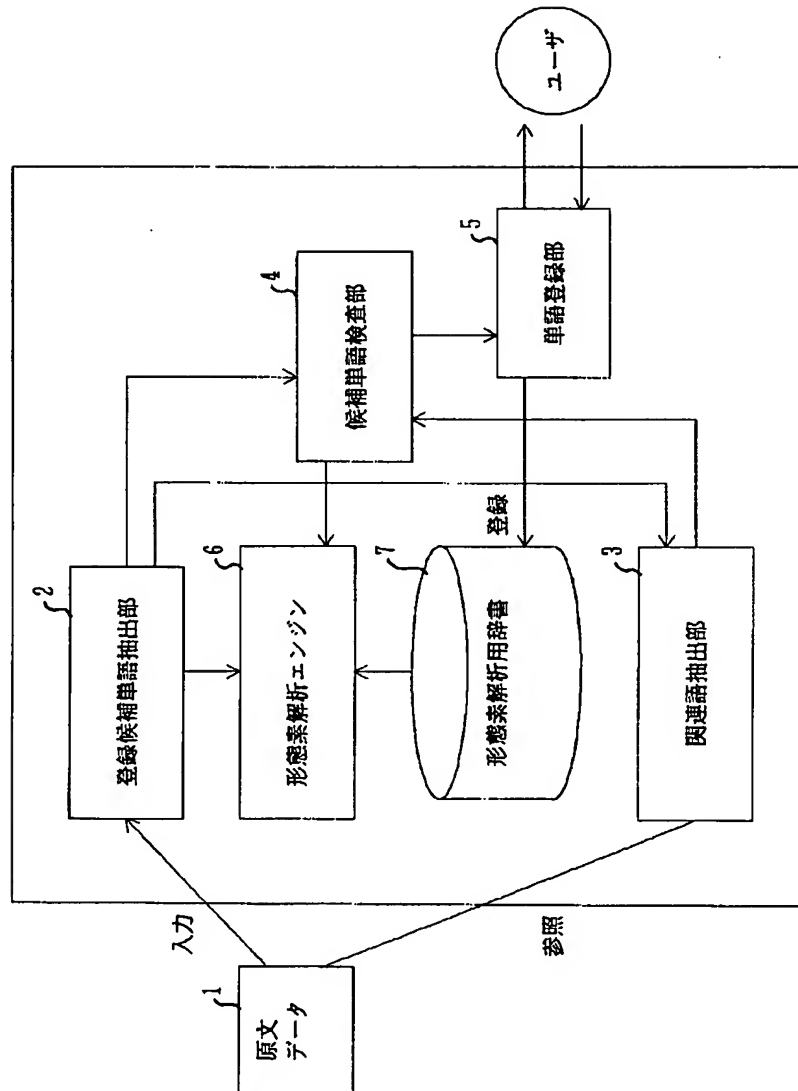
【図 1】

本発明の原理説明図



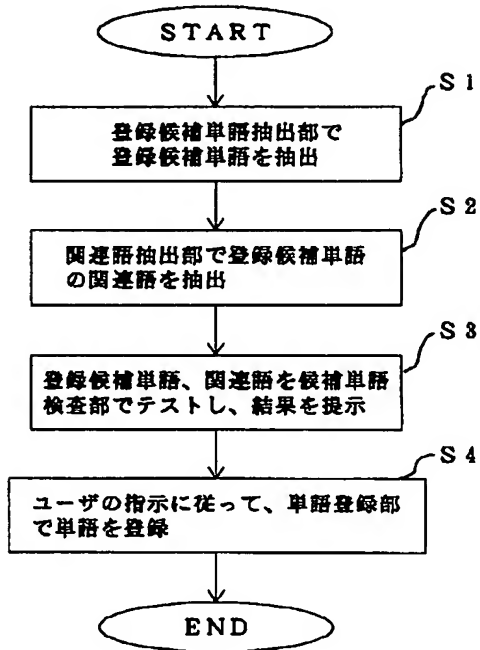
【図2】

装置構成図



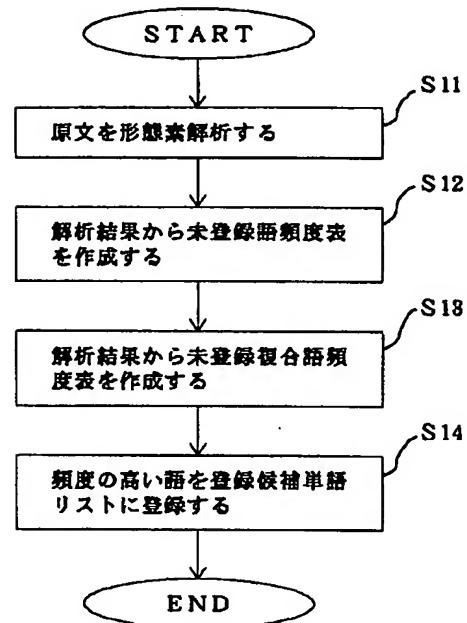
【図3】

全体の処理手順の説明図



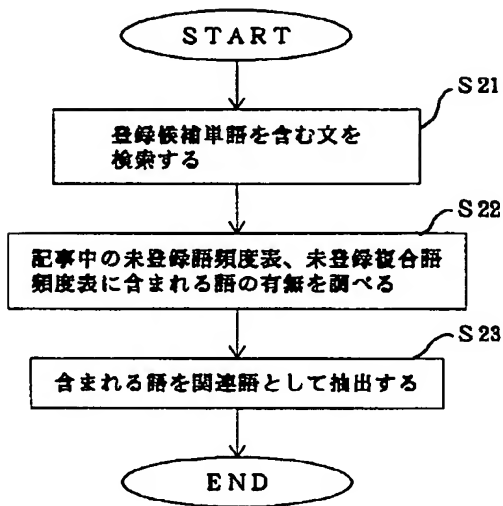
【図4】

登録候補単語抽出部の処理の説明図



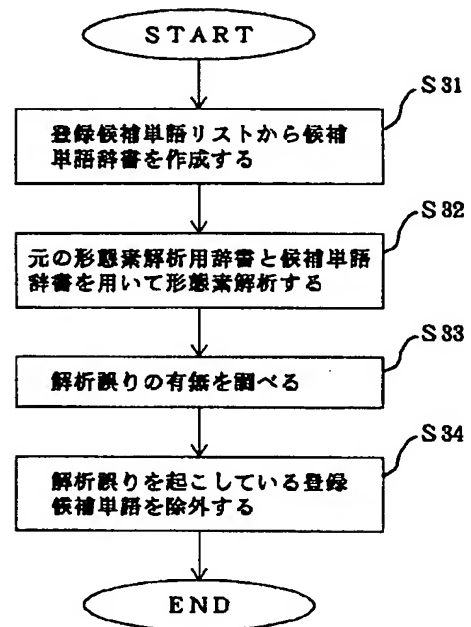
【図5】

関連語抽出部の処理の説明図



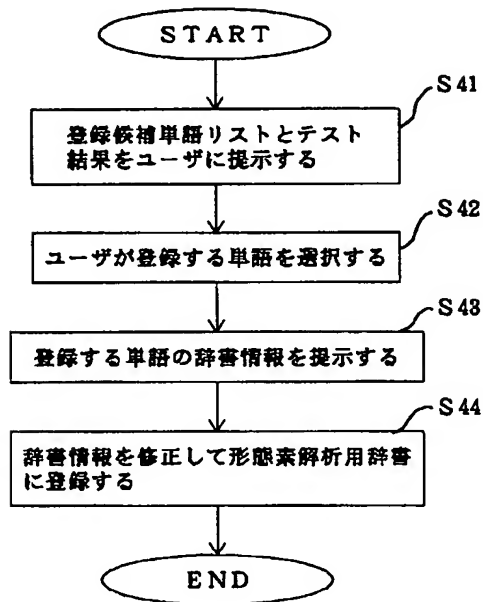
【図6】

候補単語検査部の処理の説明図



【図7】

単語登録部の処理の説明図



【図13】

関連語を登録する場合の説明図(1)

(a) 候補単語辞書(関連語)の説明

表記	品詞
園林長	普通名詞
園林長官狙撃事件	普通名詞
アウン・タン・スー・チー	普通名詞

(b) 登録前の形態素解析結果の説明

品詞	詳細品詞	表記
名詞	固有名詞	警察庁 の 園林長 官 狙撃 事件 の 捜査 を めぐ る 対応 が 適切 で ない 。
助詞	助詞「の」	
名詞	未登録語	
接尾語	接尾語	
名詞	サ変名詞	
名詞	普通名詞	
助詞	助詞「の」	
名詞	サ変名詞	
助詞	格助詞	
動詞	ラ行五段動詞	
動詞語尾	ラ行五段終止連	
名詞	サ変名詞	
助詞	格助詞	
形容動詞	ナ活用形容動詞	
形容動詞	形容動詞連用形	
助動詞	ない	
形容詞語	形容詞終止連体	
その他	句点	

【図8】

登録候補単語を登録する場合の説明図(1)

(a) 一文の形態素解析例の説明

品詞	詳細品詞	表記
名詞	サ変名詞	逮捕
動詞語尾	サ変未然形	され
助動詞	れる	れた
助動詞	た終止連体形	の
助詞	助詞「の」	は
助詞	提題助詞	オーム
名詞	普通名詞	真理
名詞	普通名詞	教
名詞	普通名詞	の
助動詞	助詞「の」	信者
名詞	人称名詞	の
助動詞	助詞「の」	林
名詞	人名	春男
名詞	人名	容疑者
名詞	人称名詞	です
助動詞	だ・です終止連	です
その他	句点	。

(b) 未登録単語頻度表の説明

頻度	品詞	詳細品詞	表記
10	名詞	未登録語	ヤンゴン
5	名詞	未登録語	園林長
2	名詞	未登録語	アコヤ

【図9】

登録候補単語を登録する場合の説明図（2）

（a）候補単語辞書の説明

表記	品詞	詳細品詞
ヤンゴン	名詞	普通名詞

（b）登録前の形態素解析結果の説明

品詞	詳細品詞	表記
名詞	地名	ミャンマー
助詞	助詞「の」	の
名詞	普通名詞	首都
名詞	未登録語	ヤンゴン
助詞	格助詞	で
名詞	人称名詞	学生
助詞	助詞「の」	の
名詞	サ変名詞	デモ
助詞	格助詞	が
動詞	ラ行五段動詞	始ま
動詞語尾	ラ行五段連用形	った
助動詞	た終止連体形	た
その他	句点	。

【図10】

登録候補単語を登録する場合の説明図(3)

(a) 「ヤンゴン」を登録した場合の形態素解析結果の説明

品詞	詳細品詞	表配
名詞	地名	ミャンマー
助詞	助詞「の」	の
名詞	普通名詞	首都
名詞	普通名詞	ヤンゴン
助詞	格助詞	で
名詞	人称名詞	学生
助詞	助詞「の」	の
名詞	サ変名詞	デモ
助詞	格助詞	が
動詞	ラ行五段動詞	が始ま
動詞語尾	ラ行五段連用形	った
助動詞	た終止連体形	った
その他	句点	。

(b) ユーザが修正した候補単語辞書の説明

表配	品詞	詳細品詞
ヤンゴン	名詞	地名

【図11】

未登録複合語頻度表を作成する場合の説明図(1)

(a) 未登録複合語頻度表の説明

頻度	表記
12	オーム／真理／教
4	園林長／官／狙撃／事件
3	第／一／波
2	アウン／・／タン／・／スー／・／チー

(b) 候補単語辞書の説明

表記	品詞	詳細品詞
オーム真理教	名詞	普通名詞

(c) 登録前の形態素解析結果の説明

品詞	詳細品詞	表記
名詞	普通名詞	オーム
名詞	普通名詞	真理
名詞	普通名詞	教
助詞	助詞「の」	の
名詞	人称名詞	信者
助詞	助詞「の」	の
名詞	人名	林
名詞	人名	春男
名詞	人称名詞	容疑者
助詞	格助詞	が
名詞	時詞	きょう
名詞	サ変名詞	逮捕
動詞語尾	サ変未然形	され
助動詞	れる	まし
助動詞	ます連用形	まし
助動詞	た終止連体形	た
その他	句点	。

【図 1 2】

未登録複合語頻度表を作成する場合の説明図 (2)

(a) 登録した後の形態素解析結果の説明

品詞	詳細品詞	表記
名詞	普通名詞	オーム真理教
助詞	助詞「の」	の
名詞	人称名詞	信者
助詞	助詞「の」	の
名詞	人名	林
名詞	人名	春男
名詞	人称名詞	容疑者
助詞	格助詞	が
名詞	時間	きょう
名詞	サ変名詞	逮捕
動詞語尾	サ変未然形	さ
助動詞	れる	れ
助動詞	ます連用形	まし
助動詞	た終止連体形	た
その他	句点	。

(b) ユーザが修正した候補単語辞書の説明

表記	品詞	詳細品詞
オーム真理教	名詞	固有名詞

【図14】

関連語を登録する場合の説明図(2)

(a) 「園林長」を登録した場合の形態素解析結果の説明

品詞	詳細品詞	表記
名詞	固有名詞	警察庁
助詞	助詞「の」	の
名詞	普通名詞	園林長
接尾語	接尾語	官
名詞	サ変名詞	狙撃
名詞	普通名詞	事件
助詞	助詞「の」	の
名詞	サ変名詞	捜査
助詞	格助詞	を
動詞	ラ行五段動詞	めぐる
動詞語尾	ラ行五段終止連	る
名詞	サ変名詞	対応
助詞	格助詞	が
形容動詞	ナ活用形容動詞	適切
形容動詞	形容動詞連用形	で
助動詞	ない	ない
形容詞語	形容詞終止連体	い
その他	句点	。

(b) 「園林長官狙撃事件」を登録した場合の形態素解析結果の説明

品詞	詳細品詞	表記
名詞	固有名詞	警察庁
助詞	助詞「の」	の
名詞	普通名詞	園林長官狙撃事件
助詞	助詞「の」	の
名詞	サ変名詞	捜査
助詞	格助詞	を
動詞	ラ行五段動詞	めぐる
動詞語尾	ラ行五段終止連	る
名詞	サ変名詞	対応
助詞	格助詞	が
形容動詞	ナ活用形容動詞	適切
形容動詞	形容動詞連用形	で
助動詞	ない	ない
形容詞語	形容詞終止連体	い
その他	句点	。

【図15】

関連語を登録する場合の説明図(3)

(a) ユーザが修正した候補単語辞書の説明

表記	品詞	詳細品詞
園林長官狙撃事件	名詞	固有名詞

(b) 登録前の形態素解析結果の説明

品詞	詳細品詞	表記
名詞	未登録語	アウン
記号	記号	・
名詞	普通名詞	タン
記号	記号	・
名詞	未登録語	スー
記号	記号	・
名詞	未登録語	チー
接尾語	人称名詞接尾語	さん
助詞	助詞「の」	の
名詞	普通名詞	勢力
助詞	格助詞	と
助詞	提題助詞	は
名詞	普通名詞	線
助詞	格助詞	を
動詞	サ変動詞	画
動詞語尾	サ変連用形	し
その他	た形接続語	て
動詞語尾	一段終止連体形	い
その他	句点	。

【図16】

関連語を登録する場合の説明図（４）

(a) 「アウン・タン・スー・チー」を登録した場合の形態素解析結果の説明

品詞	詳細品詞	表記
名詞	普通名詞	アウン・タン・スー・チー
接尾語	人称名詞接尾語	さん
助詞	助詞「の」	の
名詞	普通名詞	勢力
助詞	格助詞	と
助詞	議題助詞	は
名詞	普通名詞	一線
助詞	格助詞	を
動詞	サ変動詞	面
動詞語尾	サ変連用形	し
その他	た形接続語	てい
動詞語尾	一段終止連体形	る
その他	句点	。

(b) ユーザが修正した候補単語辞書の説明

表記	品詞	詳細品詞
アウン・タン・スー・チー	名詞	人名